1 METHOD AND DEVICE FOR PROCESSING A SPEECH SIGNAL FOR ROBUST
2 SPEECH RECOGNITION

3 The invention relates to a method and a device for processing a
4 speech signal, which is tainted by noise, for subsequent speech
5 recognition.

6 Speech recognition is being used to an increasing extent to
7 facilitate the operation of electrical devices.
8 To enable speech to be recognized what is known as an acoustic
9 model must be created. To this end, speech commands are
10 trained, a process which can be undertaken for example - for
11 the case of speaker-independent speech recognition - at the
12 factory. Training here is taken to mean the creation of so-
13 called feature vectors describing the voice command, based on
14 speaking a voice command numerous times. These feature vectors
15 (which are also called prototypes) are then collected into the
16 acoustic model, for example a so-called HMM (Hidden Markov
17 Model).
18 The acoustic model serves to determine from a given sequence of
19 speech commands or words selected from the vocabulary, the
20 likelihood of the observed feature vectors (during the
21 recognition).

22 For speech recognition or recognition of flowing speech, in
23 addition to an acoustic model a so-called speech model is also
24 used, which specifies the likelihood of individual words
25 following each other in the speech to be recognized.

26 The aim of current improvements in speech recognition is to
27 gradually achieve better speech recognition rates, i.e. to
28 increase the likelihood that a word or speech command spoken by
29 a user of the mobile communication device being recognized
30 correctly.
31 Since this speech recognition has a multiplicity of uses, it is

1 also used in environments which are adversely affected by

2 noise. In this case the speech recognition rates fall

3 drastically since the feature vectors to be found in the

4 acoustic model, for example in the HMM, have been created on

5 the basis of clean speech, i.e. speech untainted by noise. This

6 leads to unsatisfactory speech recognition in loud

7 environments, such as on the street, in busy buildings or also

8 in the car.

9 Using this prior art as its starting point, the object of the

10 invention is to create an option for also performing speech

11 recognition with a high speech recognition rate in noisy

12 environments.

13 This object is achieved by the features of the independent

14 claims. Advantageous further developments are the object of the

15 dependent claims.

16 The core of the invention is that processing of the speech

17 signal is undertaken before this signal is routed to a speech

18 recognition system for example. The speech signal undergoes

19 noise suppression within the framework of this processing.

20 Subsequently the speech signal is normalized as regards its

21 signal level. The speech signal in this case comprises one or

22 more speech commands.

23 This has the advantage that the speech recognition rates for a

24 speech command for a speech signal with noise-tainted speech

25 pre-processed in this manner are significantly higher than with

26 conventional speech recognition with noise-tainted speech

27 signals.

28 Optionally, after noise suppression, the speech signal can also

29 be fed to a unit for determining the speech activity. On the

30 basis of this noise-reduced speech signal it is then

1 established whether speech or a pause between speech is
2 present. Depending on this decision, the normalization factor
3 for signal level normalization is then determined. In
4 particular the normalization factor can be defined so that
5 pauses between speech are more heavily suppressed. Thus the
6 difference between speech signal sections in which speech is
7 present and those sections in which no speech is present
8 (pauses), becomes even more clear. This makes speech
9 recognition easier.

10 A method with the features described above can also be applied
11 to so-called distributed speech recognition systems. A
12 distributed speech recognition system is characterized by not
13 all steps within the framework of speech recognition being
14 performed in the same component. More than one component is
15 thus required. For example one component can be a communication
16 device and a further component can be an element of a
17 communication network. In this case for example the speech
18 signal detection takes place in a communication device equipped
19 as a mobile station, but the actual speech recognition on the
20 other hand is undertaken in the communication network element
21 on the network side.

22 This method can be applied both in speech recognition and also
23 when the acoustic model is being created, for example an HMM.
24 Application of the method during the creation of the acoustic
25 model in conjunction with speech recognition, based on an
26 inventively preprocessed signal, shows a further improvement in
27 the speech recognition rate.

28 Further advantages of the invention are shown with reference to
29 selected exemplary embodiments which are also illustrated in
30 the Figures.

31 The figures show:

| 1 | Fig. 1: | a histogram in which speech signals containing |
| 2 | | one or more speech commands are plotted in |
| 3 | | relation to their signal level, for the case of |
| 4 | | training to create an acoustic model; |
| 5 | Fig. 2: | a histogram of speech signals in relation to |
| 6 | | their signal level for the case of a speech |
| 7 | | recognition; |
| 8 | Fig. 3: | a schematic embodiment of an inventive |
| 9 | | processing sequence; |
| 10 | Fig. 4: | a histogram, in which the noise-reduced and |
| 11 | | speech level-normalized speech signal is |
| 12 | | plotted against the speech signal level; |
| 13 | Fig. 5: | a histogram, in which the noise-reduced speech |
| 14 | | signal is plotted against the signal level; |
| 15 | Fig. 6 | a histogram, in which the speech signal is |
| 16 | | preprocessed in the training in accordance with |
| 17 | | the invention; |
| 18 | Fig. 7 | a distributed speech processing scheme; |
| 19 | Fig. 8 | an electrical device which can be used within |
| 20 | | the framework of distributed speech processing. |

21 Fig. 8 shows an electrical device embodied as a mobile
22 telephone or mobile station MS. It has a microphone M for
23 accepting speech signals containing speech commands, a Central
24 Processing Unit CPU for processing the speech signals and a
25 radio interface FS for transmitting data, for example processed
26 speech signals.

27 The electrical device can, on its own or in combination with
28 other components, implement speech recognition with regard to
29 the accepted or detected speech commands.

30 The detailed investigations which have led to the invention
31 will now be presented:

1 Fig. 1 shows a histogram in which speech signals containing one
2 or more speech commands are sorted in respect of their signal
3 level L and this frequency H has been plotted against the
4 signal level. In this case a speech signal S, as indicated in
5 the following Figures for example, contains one or more speech
6 commands. For the sake of simplicity it is assumed below that
7 the speech signal contains a speech command. A speech command
8 can for example be formed for an electrical device equipped as
9 a mobile telephone by the request "call" as well as optionally
10 by a specific name. A speech command must be trained for speech
11 recognition, i.e. based on repeated speaking of the speech
12 command one feature vector or a number, i.e. more than one
13 feature vector is created. This training is undertaken within
14 the framework of creating the acoustic model, for example the
15 HMM, which occurs at the production stage. These feature
16 vectors are included later for speech recognition.

17 The training of speech commands which is used for the creation
18 of feature vectors is performed at a defined signal level or
19 volume level (single level training). In order to exploit the
20 dynamic range of the AD converter to convert the speech signal
21 into a digital signal, the preferred operational level is
22 around -26 dB. The definition in Decibels (dB) is produced by
23 the bits available for signal level. Thus 0 dB would mean an
24 overflow (that is exceeding the maximum volume or the maximum
25 level). Alternatively instead of a single level training,
26 training can be performed at a number of levels, for example at
27 -16, -26 and -36 dB.

28 Fig. 1 in this case shows the frequency distribution of the
29 speech level for a speech command for training.

30 A mean signal level $X_{mean}$ as well as a certain distribution of
31 the levels of the speech signal is produced for a speech
32 command. This can be represented as a Gaussian function with

1   the mean signal level $X_{mean}$ and a variance $\sigma$.

2   After the distribution of the speech commands for the training

3   situation has been seen in Fig 1, the situation for speech

4   recognition is shown in Fig 2 which again presents the

5   frequency H plotted against the signal level L in accordance

6   with Fig 1: Here the speech signal S' with one or more speech

7   commands, as is indicated in the subsequent Figures, is sorted

8   as regards its signal level L and the frequency H is plotted.

9   Because of the environmental effects, even after noise

10   reduction NR has already been applied (cf. Fig. 3) a

11   distribution shifted in relation to the training situation in

12   Fig 1 is produced, with a new mean signal level $X_{mean}$ shifted in

13   relation to the mean value $X_{mean}$ in the training.

14   It has been shown in investigations that the speech recognition

15   rate reduces drastically as a result of this shifted mean

16   signal level $X_{mean}$.

17   This can be seen from Table 1 below:

18   Table 1: Training with clean speech at different volume levels

19   or signal levels (multi-level).

20   The speech recognition rates relate to the test speech which

21   was normalized at the signal levels -16, -26, -36 dB.

| Test speech Signal levels | Word recognition rates [%] | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | Subway | | Babble | | Car | | Exhibition | |
| | Clean | 5 dB | Clean | 5 dB | Clean | 5 dB | Clean | 5 dB |
| -16 dB | 98.83 | 80.14 | 98.79 | 86.99 | 98.72 | 88.01 | 99.11 | 79.76 |
| -26 dB | 99.14 | 85.66 | 99.15 | 76.66 | 99.19 | 91.35 | 99.35 | 85.00 |
| -36 dB | 99.39 | 85.05 | 99.21 | 82.41 | 99.28 | 89.41 | 99.57 | 85.47 |

22   Table 1 lists the speech recognition rate or word recognition

23   rate for different noise environments in which training with a

24   clean speech at different volumes has been undertaken. The test

1 speech, that is the speech signal from Fig. 1, has been
2 normalized at three different levels at -16 dB, -26 dB and -36
3 dB. The speech recognition rates for different types of noises
4 with a noise level of 5 dB are shown for this different test
5 speech energy level. The different noises involved are typical
6 background noises such as subway, so-called babble noise, e.g.
7 a cafeteria environment with speech and other noises, the
8 background noise in a car as well as the noise at an exhibition
9 (i.e. similar to bubble noise, but worse, possible with
10 announcements, music etc). It can be seen from Table 1 that
11 speech recognition in noise-free speech is largely unaffected
12 by variations in the test speech energy level. However for
13 noise-tainted speech a significant reduction in speech
14 recognition can be seen. The terminal-based pre-processing
15 method AFE has been included for speech recognition here which
16 is used to create the feature vectors.

17 For the speech recognition rates investigated in Table 1 –
18 which are still not satisfactory – the situation is however
19 significantly improved compared to the speech recognition based
20 on training with only one volume level.

21 In other words the effect which an ambient noise has on an
22 acoustic model which was created on the basis of only one
23 volume of the training speech is even more plainly detrimental.

24 This has led to the inventive improvements presented below:

25 Fig. 3 now presents the execution sequence in accordance with
26 one exemplary embodiment of the invention. The speech command
27 or speech signal S, e.g. a word spoken by a person, is
28 subjected to a noise reduction NR. After this noise reduction
29 NR a noise-reduced speech signal S' is present.

30 The noise-reduced speech signal is subsequently subjected to a

1  signal level normalization SLN. This normalization is used to
2  establish a signal level which is comparable with the average
3  signal level shown in Fig. 1 by $X_{mean}$. It has been shown that
4  higher speech recognition rates can be obtained for comparable
5  mean signal levels. This means that the speech recognition rate
6  is already increased by this shifting of the signal level.

7  After the signal level normalization SLN a normalized and
8  noise-reduced speech signal S" is present. This can be
9  subsequently used for example for a speech recognition SR with
10 a higher speech recognition rate than for original test speech
11 tainted by noise.

12 Optionally the noise-reduced signal S' is split up and also
13 flows in addition to the signal level normalization SLN to a
14 Voice Activity Detection VAD unit. Depending on whether speech
15 or a speech pause is present, the normalization level with
16 which the noise-reduced speech signal was normalized, is set.
17 For example in speech pauses a smaller multiplicative
18 normalization factor can be used by which the signal level of
19 the noise-reduced speech signal S' is reduced more in speech
20 pauses than if speech is present. This means that a stronger
21 distinction between speech, that is between individual speech
22 commands for example and speech pauses is possible, which
23 further greatly improves a downstream speech recognition as
24 regards the speech recognition rate.

25 Furthermore there is provision to change the normalization
26 factor not only between speech pauses and speech sections but
27 also to vary it within a word for different speech sections.
28 The speech recognition can also be improved in this way since a
29 number of speech sections, because of the phonemes contained
30 within them, exhibit a very high signal level, for example with
31 plosive sounds (e.g. p), whereas others are rather inherently
32 silent.

1   Different methods are employed for signal level normalization,

2   for example a real-time energy normalization, as described in

3   the Article "Robust Endpoint Detection and Energy Normalization

4   for Real-Time Speech and Speaker recognithm" by Qi Li et al.

5   in IEEE Transactions on Speech and Audio Processing Vol. 10,

6   No. 3, March 2002 in Section C (P. 149-150). A further signal

7   level normalization method is described within the framework of

8   the ITU, which can be found under ITU-T, ''SVP56: The Speech

9   Voltmeter'', in software Tool Library 2000 User's Manual, pages

10   151-161, Geneva. Switzerland, December 2000. The normalization

11   described in this document works "off-line" or in what is known

12   as "batch mode", i.e. not simultaneously or contemporaneously

13   with speech recognition.

14   For noise reduction NR (cf. Fig. 3) different known methods are

15   also provided, for example methods operating in the frequency

16   area One such method is described in "Computationally efficient

17   speech enhancement using RLS and psycho-acoustic motivated

18   algorithm" by Ch. Beaugeant et al. in Proceedings of 6th World

19   Multi-conference on Systemics, Cybernetics and Informatics,

20   Orlando 2002. The system described in this document is based on

21   an analysis-by-synthesis system in which the parameters

22   describing the (clean) speech signal and the noise signal are

23   extracted frame-by frame recursively (cf. Section 2 "Noise

24   Reduction in the Frequency Domain", Section 3 "Recursive

25   implementation of the least square algorithm" in this

26   document). The clean speech signal thus obtained is further

27   weighted (cf. Section 4 "Practical RLS Weighting Rule") and an

28   estimation of the power of the noise signal is undertaken (cf.

29   Section 5 "Noise Power Estimation"). Optionally the results

30   obtained can be refined by means of psychoacoustic motivated

31   methods (Section 6:"Psychoacoustic motivated method"). Further

32   noise reduction methods which can be included in accordance

33   with an embodiment shown in Fig. 3 are for example described in

1 ETSI ES 202 0505 V1.1.1 dated October 2002 in Section 5.1
2 ("Noise Reduction").

3 An unprocessed speech signal S as regards noise reduction NR
4 and signal normalization is used as the basis for the frequency
5 distributions in Fig. 1 (training situation) and 2 (test
6 situation, i.e. for a speech recognition). The noise-reduced
7 speech signal S' is used as a basis for the frequency
8 distribution in Figure 5. The noise-reduced and signal-level-
9 normalized signal is used as the basis for the distributions in
10 Figures 4 (test situation) and 5 (training situation).

11 The idea underlying the schematic execution sequence shown in
12 Fig 3 of a speech signal processing for a subsequent speech
13 recognition is presented in Figures 4 to 6.

14 Fig. 5 shows a frequency distribution for a noise-reduced
15 speech signal S' as occurs for example in Fig. 3 after the
16 noise reduction NR. Compared to Fig. 2, which relates for
17 example to the frequency distribution for a speech signal S
18 shown in Fig. 3, a further noise reduction NR has thus been
19 undertaken.

20 The center of the frequency distribution of this noise-reduced
21 speech signal S' compared to the speech level L is to be found
22 at a mean level $X_{mean}$. The distribution has a width $\sigma'$. In the
23 transition to Fig. 4, signal normalization SLN is performed on
24 the noise-reduced signal S' shown in Fig. 5. This means for
25 example that the speech signal used as a basis for the
26 distribution in Fig. 4 would correspond to the noise-reduced
27 and signal-level-normalized speech signal S".

28 A signal level normalization brings the actual signal level in
29 Fig. 5, to a desired signal level, for example the signal level
30 obtained in training, indicated in Fig. 1 by $X_{mean}$. Furthermore

1   signal level normalization SLN leads to the distribution

2   becoming narrower i.e. to $\sigma''$ being narrower than $\sigma'$. This means

3   that the mean signal level $X_{mean}''$ in Fig. 4 can more easily be

4   reconciled with the mean signal level $X_{mean}$ in Fig. 1, which was

5   obtained in training. This leads to higher speech recognition

6   rates.

7   The application of what has been explained above is now

8   examined for speech recognition in conjunction with Fig. 7.

9   As already explained at the start, the speech recognition can

10   take place in one component or distributed amongst a number of

11   components.

12   For example in an electrical device which is embodied as a

13   mobile station MS there can been means for recognizing the

14   speech signal, e.g. the microphone M shown in Fig. 8, means for

15   a noise reduction NR and means the signal normalization SN. The

16   latter can be implemented within the framework of the Central

17   Processing Unit CPU. Thus the idea presented in Fig. 3 of

18   speech signal processing in accordance with one embodiment of

19   the invention as well as the subsequent speech recognition in a

20   mobile station can be implemented on its own or in conjunction

21   with an element of a communication network.

22   In accordance with an alternative embodiment the speech

23   recognition SR (see Fig. 3) is even undertaken on the network

24   side. To this end the feature vectors created from a speech

25   signal S" are transmitted via a channel, especially a radio

26   channel to a central unit in the network. Here the speech

27   recognition is undertaken on the basis of the transmitted

28   feature vectors especially on the basis of the model created

29   during production. During production can mean especially that

30   the acoustic model is created by the network operator.

31   In particular the proposed speech recognition can be applied to

1 speaker-Independent speech recognition as is for example

2 undertaken within the framework of the so-called Aurora

3 scenario.

4 A further improvement emerges in speech commands are already

5 normalized when the acoustic model is created during production

6 or during training in respect of the signal level. This means

7 that the distribution of the signal level is namely narrower,

8 by which an even better match between the distribution shown in

9 Fig. 4 and the distribution achieved in the training is

10 obtained. Such a distribution of the frequency H in relation to

11 the signal level L for a speech command in training for which

12 signal level normalization has already been performed, is shown

13 in Fig. 6. The mean training level $X_{mean\_new}$ last produced

14 matches the mean level $X_{mean}''$ (Fig. 4) of the noise-reduced and

15 signal-level-normalized speech signal S" (Fig. 3). As has

16 already been shown, a match in the mean levels is one of the

17 criteria for a high speech recognition rate. Furthermore the

18 width of the distribution in Fig. 6 is very narrow which makes

19 it easier to reconcile this distribution with the distribution

20 in Fig. 4, i.e. bring it to the same signal level.

21 Fig. 7 shows a Distributed Speech Recognition (DSR). A

22 distributed speech recognition can for example be used within

23 the framework of the AURORA project of the ETSI STQ (Speech

24 Transmission Quality) already mentioned.

25 With a distributed speech recognition a speech signal, for

26 example a speech command, is detected at a unit and feature

27 vectors describing this speech signal are created. These

28 feature vectors are transmitted to another unit, typically a

29 network server. Here the feature vectors are processed and

30 speech recognition is performed on the basis of these feature

31 vectors.

32 Fig. 7 shows a mobile station MS as a first unit or component

1 and a network element NE.

2 The mobile station MS, which is also referred to as a terminal,
3 features means AFE for terminal-based preprocessing which are
4 used to create the feature vectors. For example the mobile
5 station MS is a mobile radio device, portable computer or any
6 other mobile communication device. The means AFE for terminal-
7 based preprocessing is for example the "Advanced Front End"
8 discussed within the framework of the AURORA project.

9 The means AFE for terminal-based preprocessing comprises means
10 for standard processing of speech signals. This standard speech
11 processing is for example described in Specification ETSI ES
12 202050 V1.1.1 dated October 2002 in Fig. 4.1. On the mobile
13 station side the standard speech processing includes feature
14 extraction with the steps noise reduction, waveform processing,
15 cepstrum calculation as well as blind equalization. Feature
16 compression and preparation for transmission are subsequently
17 undertaken. This processing is known to the person skilled in
18 the art, for which reason it is not discussed in further detail
19 here.

20 In accordance with an embodiment of the invention the means AFE
21 for terminal-based preprocessing also comprises means for
22 signal level normalization and voice activity detection in
23 accordance with Fig. 3.

24 These means can be integrated into the AFE means or
25 alternatively implemented as a separate component.

26 Using subsequent means FC for feature compression, terminal-
27 based preprocessing AFE, the one or more feature vectors which
28 are created from the speech command are compressed to allow
29 them to be transmitted via a channel CH.

30 The other unit is for example formed by a network server as

1 network element NE. In this network element NS the feature

2 vectors are decompressed again using means FDC for feature

3 vector decompression. In addition means SSP are used for

4 server-side preprocessing, so that the means SR for speech

5 recognition can then be used to perform speech recognition

6 based on a Hidden Markov Model HMM.

7 The results of inventive improvements will now be explained:

8 Speech recognition rates for different training of the speech

9 commands as well as different speech levels or volumes which

10 are included for speech recognition (test speech) are shown in

11 Tables 1 to 2.

12 Table 2 now shows the speech recognition rates for different

13 energy levels of the test speech. The training is undertaken at

14 a speech energy level of -26 dB. The test speech has been

15 subjected to noise suppression and speech level normalization

16 in accordance with Fig. 3. It can be seen from Table 2 that the

17 speech recognition rates for clean speech are again

18 consistently high. The significant improvement compared to the

19 previous speech recognition method lies in the fact that the

20 difference which can be seen in Table 1 in the speech

21 recognition rates for noise-tainted speech (for a signal-to-

22 noise ratio" of 5 dB) is raised depending on the energy level

23 of the test speech. The "Advanced Front End" described above is

24 employed for speech recognition.

25 Table 2:

| Test Speech Energy levels | Word Recognition Rates [%] | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | Subway | | Babble | | Car | | Exhibition | |
| | Clean | 5 dB | Clean | 5 dB | Clean | 5 dB | Clean | 5 dB |
| -16 dB | 99.45 | 83.79 | 98.85 | 75.63 | 99.02 | 86.34 | 99.35 | 79.67 |
| -26 dB | 99.20 | 84.71 | 98.88 | 74.37 | 99.05 | 87.89 | 99.32 | 80.56 |
| -36 dB | 98.86 | 84.71 | 98.70 | 75.00 | 98.78 | 87.77 | 99.01 | 80.47 |

26

1